# Pathema

## TIGR/BRC

Owen White
Oct, 13th 2004

# Primary Focus

- Annotation
  - Gene features/GO assignments
  - Biochemical pathways
- Comparative analysis
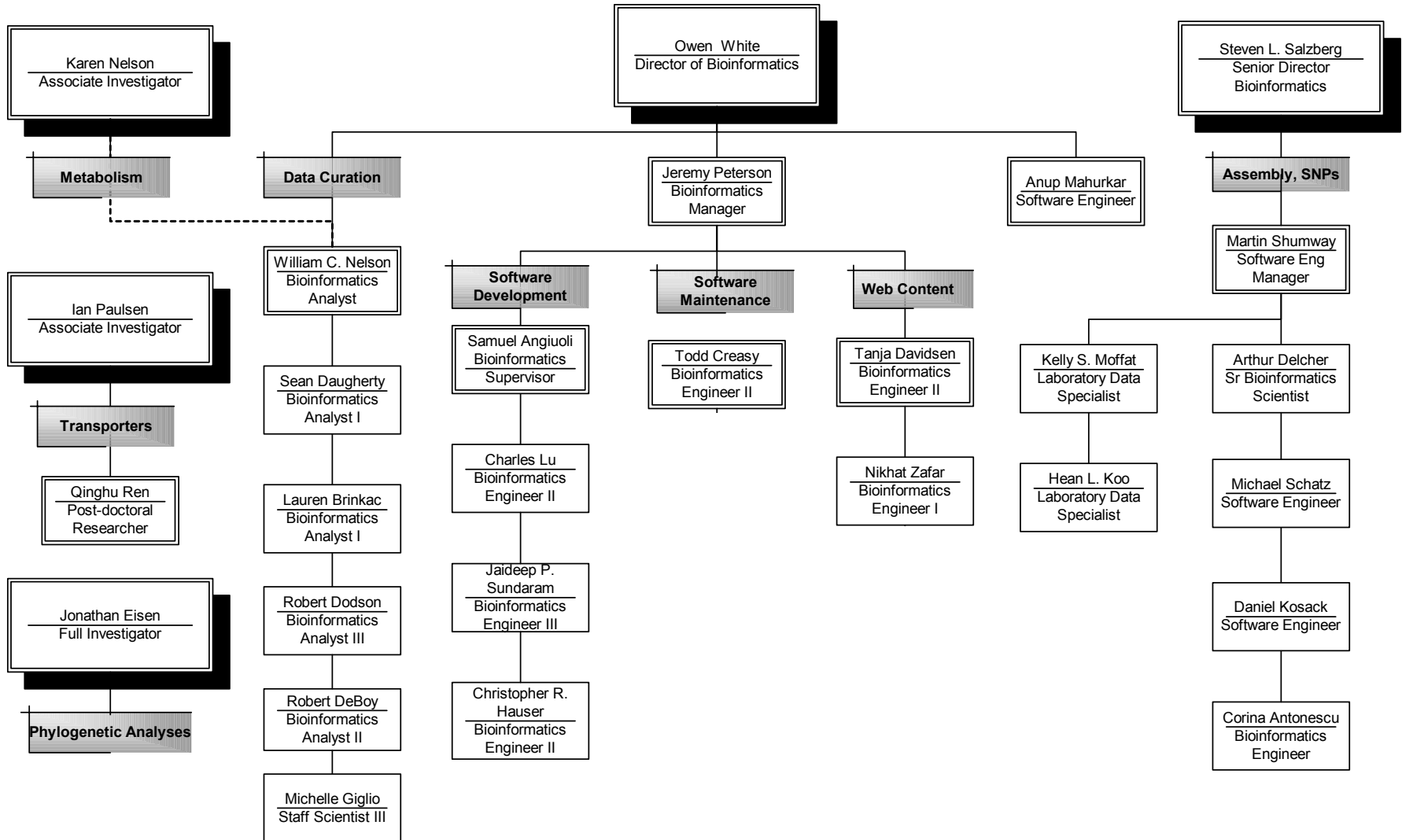- SNPs across all category A-C orgs
- Instructional Classes
  - Annotation/SOPs
  - Software
  - Web usage
- Open Source Software Support

# TIGR/BRC Organisms

- *Bacillus anthracis,*
- *Burkholderia mallei*
- *Burkholderia pseudomallei*
- *Clostridium botulinum*
- *Clostridium perfringens*
- *Francisella tularensis.*

# People

Karen Nelson
Associate Investigator

Owen White
Director of Bioinformatics

Steven L. Salzberg
Senior Director
Bioinformatics

**Metabolism**

**Data Curation**

Jeremy Peterson
Bioinformatics
Manager

Anup Mahurkar
Software Engineer

**Assembly, SNPs**

William C. Nelson
Bioinformatics
Analyst

**Software Development**

**Software Maintenance**

**Web Content**

Martin Shumway
Software Eng
Manager

Ian Paulsen
Associate Investigator

Sean Daugherty
Bioinformatics
Analyst I

Samuel Angiuoli
Bioinformatics
Supervisor

Todd Creasy
Bioinformatics
Engineer II

Tanja Davidsen
Bioinformatics
Engineer II

Kelly S. Moffat
Laboratory Data
Specialist

Arthur Delcher
Sr Bioinformatics
Scientist

**Transporters**

Lauren Brinkac
Bioinformatics
Analyst I

Charles Lu
Bioinformatics
Engineer II

Nikhat Zafar
Bioinformatics
Engineer I

Hean L. Koo
Laboratory Data
Specialist

Michael Schatz
Software Engineer

Qinghu Ren
Post-doctoral
Researcher

Robert Dodson
Bioinformatics
Analyst III

Jaideep P.
Sundaram
Bioinformatics
Engineer III

Daniel Kosack
Software Engineer

Jonathan Eisen
Full Investigator

Robert DeBoy
Bioinformatics
Analyst II

Christopher R.
Hauser
Bioinformatics
Engineer II

Corina Antonescu
Bioinformatics
Engineer

**Phylogenetic Analyses**

Michelle Giglio
Staff Scientist III

+2

# Microbial Software & Services

www.tigr.org/software

# Software

- Mummer - whole genome alignment.
- Glimmer - gene finding system.
- Manatee - manual annotation tool.
- Workflow - custom pipelines.
- Sybil - comparative analysis system.

# Mummer 3.0

- MUMs: Maximal Unique Matches
  - Algorithm finds all matches
  - String them together and align gaps
- Suffix trees
  - Fast alignment of long DNA sequences
  - Linear time and space requirements
  - Streaming algorithm
- Memory maximization
  - 2 year dev time optimizing suffix tree impl.

# Mummer Performance

On a 2.4 GHz Pentium PC running Linux:

| | 3.0 | | 2.0 | |
|---|---|---|---|---|
| | Mb | Time | Mb | Time |
| E. coli K12 vs. E. coli O157 | 78 | 13s | 102 | 14s |
| P. falciparum all chromosomes vs. P. yoelii | 552 | 16:15 | 752 | 18:29 |
| D. melanogaster arm 2L vs. D. pseudoobscura all contigs | 467 | 13:55 | 465 | 14:39 |

# TIGRFams

- Heavily curated multiple alignments based on protein families of the same function.
- Proposed "cure" for transitive annotation.
- Based on Hidden Markov Models (HMMs).
- >1,000 families.
- Complete assignments to GO.
- Cutoff scores for each family.
  - Trusted (automated name assignment)
  - Noise (manual inspection required)
- Downloadable. Fully integrated into Interpro.

# TIGRFAMs: **Orthologous Families**



|  | M. j. | Chl tr | Strep | Chl tr | Therm | A. f. | B. b. | Caul | D. r. | H. i. | H.p. | M. g. | TB | Neis. | Trep | Vib. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| TIGRFAMS | 174 | 190 | 247 | 0 | 256 | 216 | 189 | 332 | 302 | 386 | 298 | 148 | 338 | 370 | 186 | 450 |
| Identified genes | 522 | 393 | 1894 | 381 | 768 | 899 | 559 | 1338 | 947 | 1024 | 889 | 263 | 1210 | 961 | 504 | 1676 |

# 3-tier Architecture

Client X

Client Y

## API

Shared

### Application specific API

Client X

Client Y

## Database layer (SQL)

TIGR schemas

Sybase

Mysql

Chado schema

Sybase

Postgres

# 3-tier Architecture

- API
  – Unified access to data for client apps
  – Decouples client applications from SQL
- Portable
  – Cross schema/vendor support
  – Perl
- Extensible
  – Application specific API and database layer
- Thin/Lightweight
  – Single layer between client and SQL

# Manatee

- MANual Annotation Tool Etc, Etc…

# _Listeria monocytogenes 4b_ | Genome Summary

□ The Genome Summary page gives a breakdown of the characteristics of the genome. All features such as ribosomal binding sites,RNAs,phages,inteins, and terminators are shown. The page also gives a breakdown of the start sites and their frequency.Each molecule in the project is characterized showing the length, GC content. base frequencies, Percent coding and a link to the oligomernucleotide skew table.

| _Features_ | Count | feat_type |
|---|---|---|
| ▸ ribosome binding site | 2963 | RBS |
| ▸ Open Reading Frame | 2958 | ORF |
| ▸ rho-independent terminator | 752 | TERM |
| ▸ transfer RNA | 67 | tRNA |
| ▸ ribosomal RNA | 18 | rRNA |
| ▸ Bacteriophage | 2 | PHAGE |
| ▸ structural RNA | 2 | sRNA |

| _Start Sites_ | Number | Percent |
|---|---|---|
| ▸ATG: | 2425 (2292) | 82.1% (84.7%) |
| ▸GTG: | 265 (216) | 9.0% (8.0%) |
| ▸TTG: | 265 (197) | 9.0% (7.3%) |
| ▸OTHER: | 0 0 | 0.0% (0.0%) |

Numbers in parentheses do not include hypothetical proteins

## 'annot_938' Information Table

| | |
|---|---|
| ▸Assembly ID: | 942 |
| ▸Type: | chromosome |
| ▸Molecule Length: | 2905309 bp |
| ▸GC Content: | 38% |
| ▸Base Frequencies: | **(A)**   **(C)**   **(G)**   **(T)**<br>31.1%   19.1%   18.9%   30.9% |
| ▸Funny Characters: | **a**<br>1 |
| ▸Number of ORFs: | 2958 |
| ▸Average Gene Length: | 875 nt |
| ▸Percent Coding: | 89.2% |
| ▸Percent Coding OR tRNA, rRNA, or Repeat: | 89.2% |
| ▸Skew Table | |

Name: **cell growth and/or maintenance**

Type: **process**

Definition: **Any process required for the survival and growth of a cell.**

Comment: **UNDEFINED**

Synonym: **NONE**

Secondary ID: **NONE**

**Absolute Path in GO Tree: 1 instance detected**

+Ontology (TI:0000001)[R]1739
 +Gene_Ontology (GO:0003673)[P]1739
   +biological_process (GO:0008150)[P]1730
      +**cell growth and/or maintenance (GO:0008151)**[

# View Mode: Regular

+Ontology (TI:0000001)[R]1739
    +Gene_Ontology (GO:0003673)[P]1739
        +biological_process (GO:0008150)[P]1730
            +cell growth and/or maintenance (GO:0008151)[I]1316
                +transport (GO:0006810)[I]300
                +cell proliferation (GO:0008283)[I]
                +autophagy (GO:0006914)[I]
                +stress response (GO:0006950)[I]19
                 chemi-mechanical coupling (GO:0006943)[I]
                +cell motility (GO:0006928)[I]29
                +membrane fusion (GO:0006944)[I]
                +cell-cell fusion (GO:0006947)[I]
                +budding (GO:0007114)[I]
                +sporulation (GO:0030435)[I]3
                +homeostasis (GO:0019725)[I]13
                +cell organization and biogenesis (GO:0016043)[I]50
                +cell cycle (GO:0007049)[I]33
                +cell growth (GO:0016049)[I]
                +metabolism (GO:0008152)[I]930
                +regulation of cell shape and cell size (GO:0007148)[I]5
            +death (GO:0016265)[I]1
             biological_process unknown (GO:0000004)[I]380
            +viral life cycle (GO:0016032)[I]
            +physiological processes (GO:0007582)[I]14
            +development (GO:0007275)[I]18
            +cell communication (GO:0007154)[I]71

# Gene Information Page

*(overlapping titles: Gene Identification / Gene Ontology and TIGR Roles / Graphical Display of Analyses / Search PFAM for HMM role)*

**Brucella suis 1330** | **Gene Curation Page**

Help text goes here

## GENE CURATION INFORMATION

**ORFA01956 (BRA1080)**

View BER Searches
asmbl_id: 2468
Reload Page

end5/end3: 1063328 / 1062486
gene length: 843
protein length: 281
mol. wt.: 30372.93

database: gbr
feat_name/locus
New Gene

Select Function

Refresh Searches

## GENE IDENTIFICATION

Gene Name: peptide ABC transporter, permease protein

Gene Symbol

EC Number

comment: Start confidence Low. 5 GES regions. Appears to be located in a dipeptide transport operon - LMB

pub_comment:

auto_comment

## GENE ONTOLOGY
None Assigned

Add go_id | Ev_code ISS | Reference

## TIGR ROLES
142: Transport and binding proteins, Amino acids, peptides and amines

Add role_ids (separate with spaces): | Delete role_ids (click on ids above):

## -: SUBMIT DATA :-
☑ Start site edit:lbrinkac
Start confidence is high.
☑ Completed: lbrinkac

## -: EVIDENCE PICTURE :-

## HMM

PF00528: Binding-protein-dependent transport systems inner membrane component

---

## HMM

PF00528: Binding-protein-dependent transport systems inner membrane component

gene_sym: none
ec#: none
role_id: none

Isology: domain | Total score: 43.3 | Trusted cutoff: 10.00 | Noise cutoff: 9.90 | Total expect: 5.4e-09

View Alignment | Coords | HMM Coords | Score | Expect | Curation | [Add To GO Evidence]
align page | 167-243 | 1-77 / 77 | 43.3 | 5.4e-09 | ☑

## PROSITE

PS00402: Binding-protein-dependent transport systems inner membrane comp. sign.

Match sequence: LEVREREHVEAAIAAGAGSGRILFKHILP

Coords | Precision | Recall | Curation
168/196 | 0.73 | 0.63 | ☑ | [Add To GO Evidence]

## INTERPRO DATA
Query Sequence ORFA01956 - Length 281 aa.
IPR000515 PF00528 — not stored Binding-protein-dependent transport systems inner membrane
IPR001991 PR00173 — not stored Sodium:dicarboxylate

InterPro
IPR000515 Binding-protein-dependent transport systems inner membrane component
IPR001991 Sodium:dicarboxylate symporter family

## -: ATTRIBUTES :-
No Frameshifts Detected.
**View Paralogous Families**
Domain : fam_PF00528
Domain : fam_187

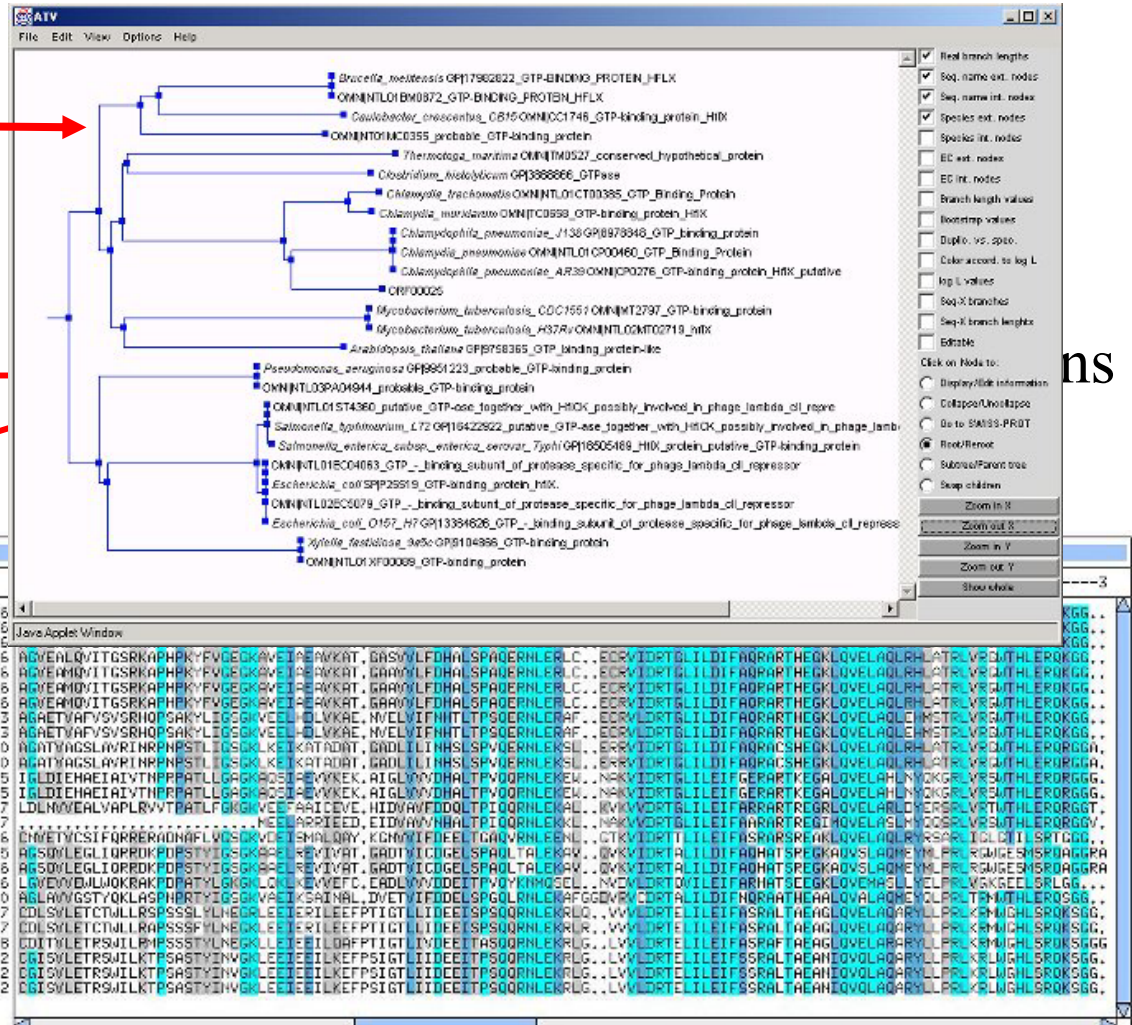## SIGNAL_P
No signalp information available. [Run signalp]

## LIPOPROTEIN INFORMATION
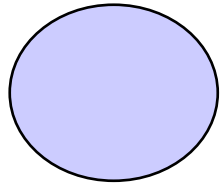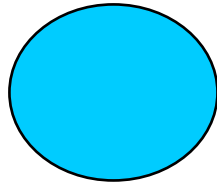No LP Information Available

# Pair-wise Alignment Summary

# Manatee: Implementation

- ~15 reports
- Linux OS
- Fully database independent.
  - Mysql
  - Sybase
- Perl Application programmer interface (API).
- Web compatible.
- Documented
  - Example databases available.
  - Online user docs.
  - Programmer/interface.
  - Installation.
- Installed at several sites/contributions from other developers received

# Chado: Open Source Database

Sequence    Comp analysis    General    Organism    CV

- Collaborative relational database ~15 people.
  - Flybase/Harvard
  - GMOD Consortium
- Composed of several modules
- Freely available as open source
- Many support tools under development by several laboratories. (see www.gmod.org)

# Annotation attributes

Annot. Attribute terms

Gene name
EC#
Gene sym
GO
…

protein

CDS

transcript

exon

gene

Genomic Contig

featureprop:
type_id =
value =

cvterm:
cv_id
cvterm_id
name

# Scaffold mappings

Feature (Gene,match)

Featureloc

featureloc:
   nbeg = z
   nend = w
   rank = 1
   locgroup = 0
   srcfeature_id = *

featureloc:
   nbeg = x
   nend = y
   rank = 0
   locgroup = 0
   srcfeature_id = *

featureloc:
   nbeg = x
   nend = y
   rank = 0
   locgroup = 1
   srcfeature_id = *

NNNN

NNNN

Contigs

Genomic scaffold

# Multiple alignments

member sequence     member sequence

Multiple
alignment
feature

featureloc

featureloc:
  nbeg = z
  nend = w
  rank = 0
  locgroup = 0
  residues = [text]
  srcfeature_id = *

featureloc:
  nbeg = z
  nend = w
  rank = 1
  locgroup = 0
  residues = [text]
  srcfeature_id = *

featureloc:
  nbeg = x
  nend = y
  rank = 2
  locgroup = 0
  residues = [text]
  srcfeature_id = *

featureloc:
  nbeg = x
  nend = y
  rank = 3
  locgroup = 0
  residues = [text]
  srcfeature_id = *

member sequence     member sequence

# Sequence variations
## SNPs (redundant mapping to protein)

SNP

featureloc

protein

I => T

featureloc:
  residue_info = "I"
  nbeg = z
  nend = z
  rank = 0
  locgroup = 1
  srcfeature_id = *

featureloc:
  residue_info = "T"
  rank = 1
  locgroup = 1

featureloc:
  residue_info = "A"
  nbeg = x
  nend = x
  rank = 0
  locgroup = 0
  srcfeature_id = *

featureloc:
  residue_info = "G"
  rank = 1
  locgroup = 0

A => G

Genomic Contig

# Chado ←→BSML

- No problem.
- Heavily documented
- Scaffolds
- Attribution
- CV terms
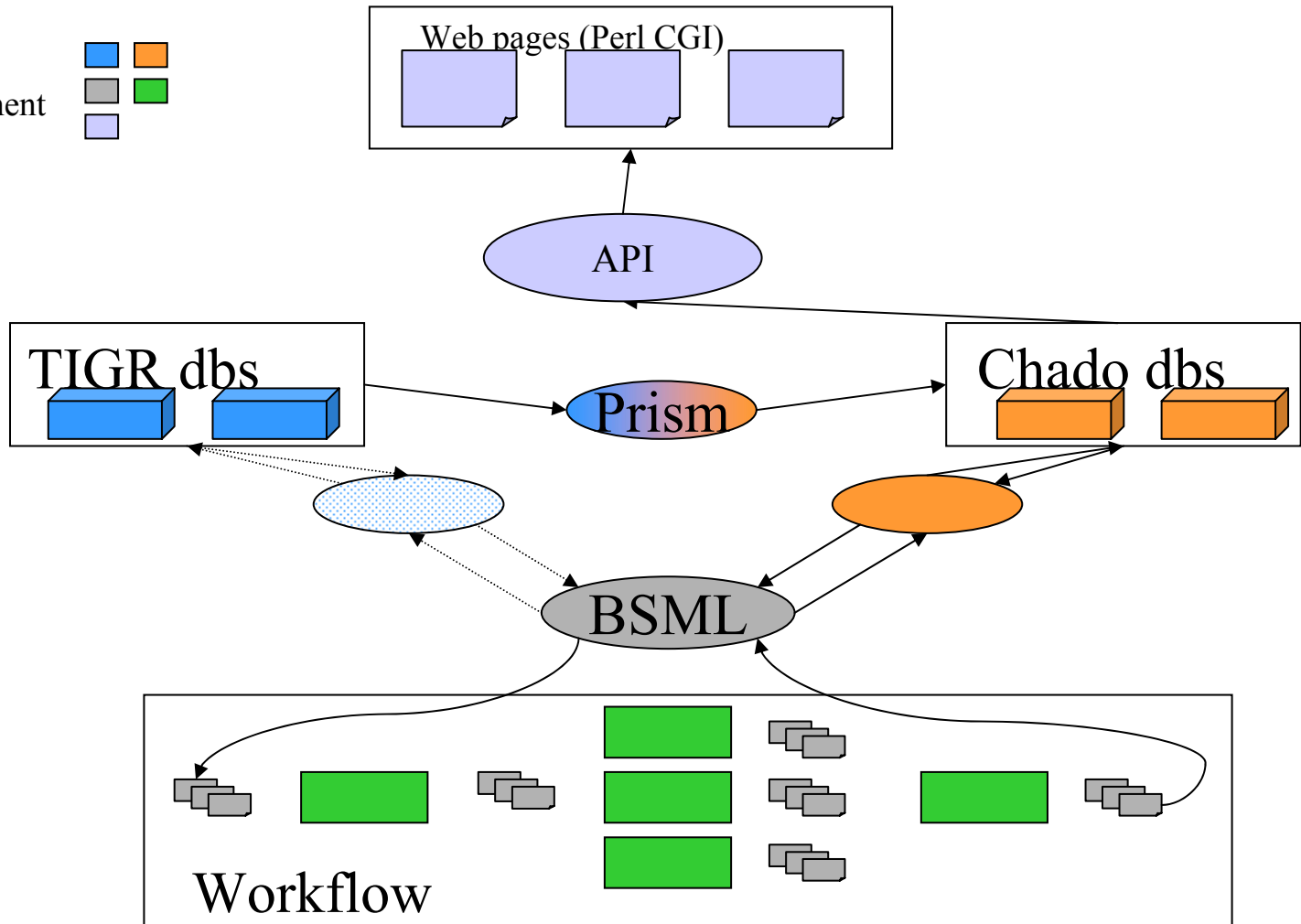- Linkage into other documents
- Large-scale file management

# Project architecture

*Database
*Data management
*Web

Web pages (Perl CGI)

API

TIGR dbs

Prism

Chado dbs

BSML

Workflow

# Workflow System

Viewing systems exist (e.g., gbrowse, Apollo, Manatee, Artemis), but how to *create* data?

- Implemented in Java, ~5000 lines
- Describes workflow as a directed acyclic graph
- Supports serial and parallel processes
- Executes onto Condor/LSF
- Two main files:
  - Human-readable config files
  - XML templates
- config file + XML template → XML "instance"
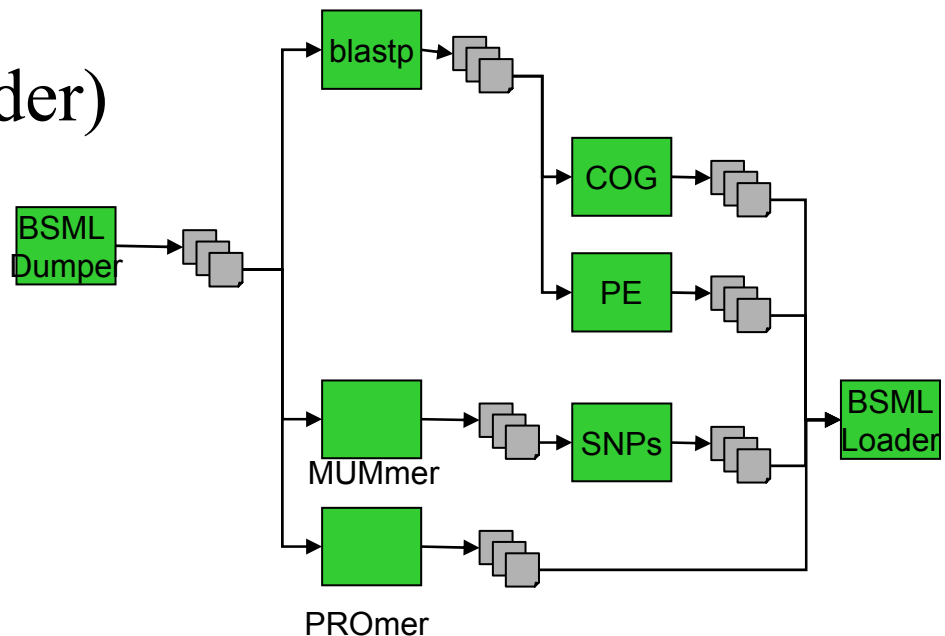
# Workflow System

- XML instance, contents:
  - Complete description of pipeline
  - Contains status of pipeline
  - Allows monitoring:
    - resumption of failed instances
    - straightforward tracking of multiple instances
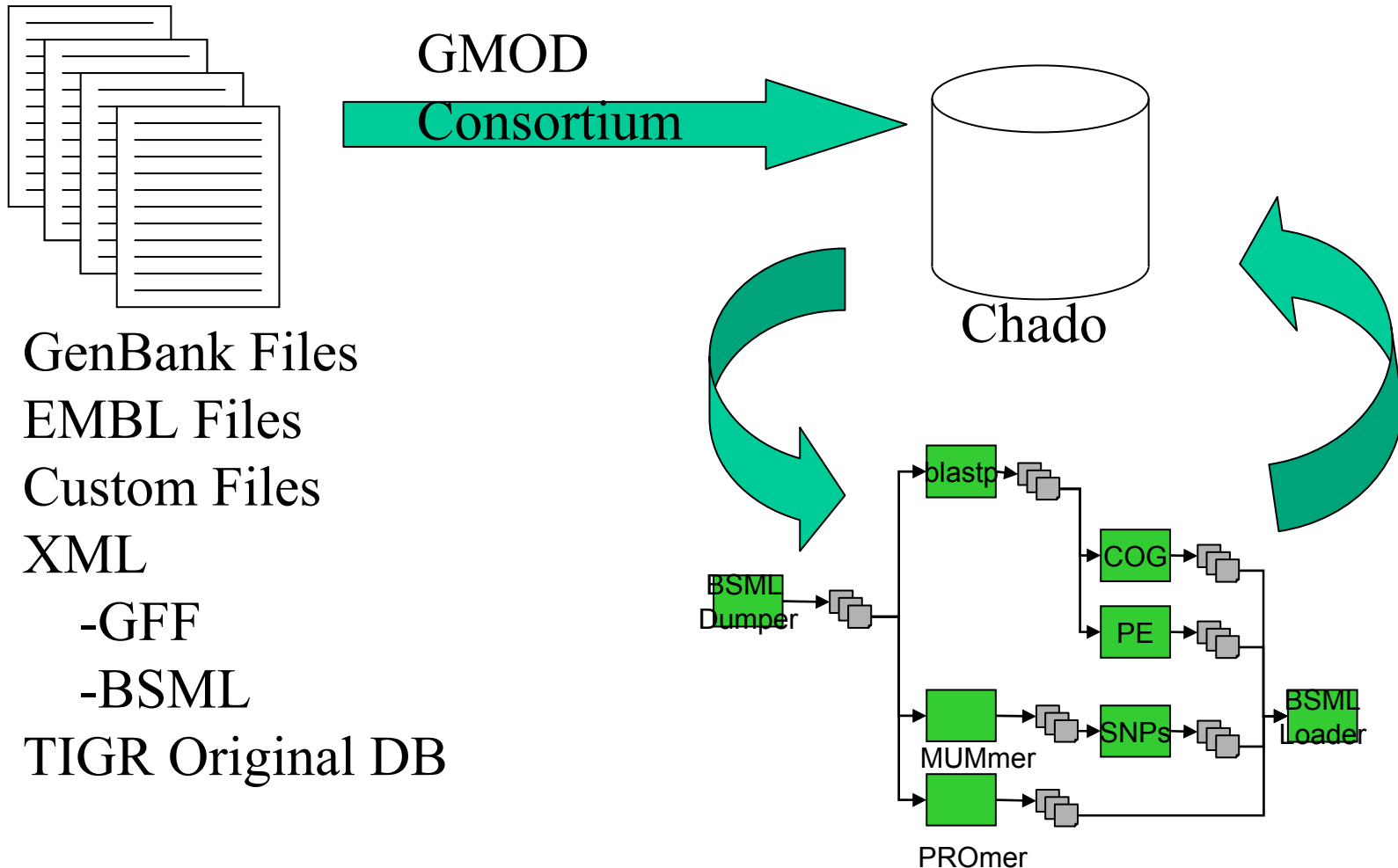
# Workflow Computes

- Blast
- Position effect (conserved gene order)
- MUMmer
  - SNPs
- PROmer
- Gene families
  - COGs
  - Paralogs



Primary output: BSML-XML

# Data Prep For Comparative Analysis

GenBank Files
EMBL Files
Custom Files
XML
  -GFF
  -BSML
TIGR Original DB

GMOD
Consortium

Chado

blastp

BSML
Dumper

COG

PE

SNPs

MUMmer

PROmer

BSML
Loader

**A  B  C**

# TIGR
## THE INSTITUTE FOR GENOMIC RESEARCH

### Microbial Synteny Tools: Match Table Display Launcher
*Streptococcus pneumoniae strains comparison* (pneumo)

**Select Match Analysis:** Position Effect ▾

**Select Molecules for Analysis** (Top molecule will be made reference molecule)

| | |
|---|---|
| S.pneumo TIGR4 chromosome<br>S.pneumo R6 chromosome<br>S.pneumo G54 pseudochromo<br>S.pyogenes M1 chromosome<br>S.agalactiae main chromos<br>S.pneumoniae 670<br>S.agalactiae h36b chromos | S.pneumo TIGR4 chromosome<br>S.pneumo R6 chromosome<br>S.pneumo G54 pseudochromo |

Add >>
<< Remove

**Click on Submit to Launch Match Analysis Tool**

Submit   Reset

**Questions?  Comments?**  Please feel free to send us [feedback](#)!

| | | | | | |
|---|---|---|---|---|---|
| | | | | ORFB01863 | hypothetical protein |
| | | | | ORFB01865 | hypothetical protein |
| ORF02100 | conserved hypothetical protein, degenerate | | | | |
| ORF02102 | xanthine phosphoribosyltransferase | NTORFA1660 | Xanthine phosphoribosyltransferase | ORFB01867 | xanthine phosphoribosyltransferase |
| ORF02103 | xanthine permease | NTORFA1661 | Nucleobase:cation symporter for xanthine | ORFB01868 | xanthine/uracil permease family protein |
| | | | | ORFB01869 | restriction endonuclease SsuRA |
| | | | | ORFB01870 | dpnA protein |
| | | | | ORFB01871 | DNA adenine methylase |
| ORF02104 | DpnD protein | NTORFA1662 | Restriction system of S. pneumoniae | | |
| ORF02105 | type II restriction endonuclease DpnI | NTORFA1663 | Type II restriction enzyme DpnI (dpnC) | | |
| ORF02107 | conserved hypothetical protein | NTORFA1664 | Conserved hypothetical protein | ORFB01872 | uncharacterized domain 1, putative |
| ORF02110 | galactose-1-phosphate uridylyltransferase | NTORFA1665 | Galactose-1-phosphate uridylyltransferase | ORFB01873 | galactose-1-phosphate uridylyltransferase |
| ORF02111 | galactokinase | NTORFA1666 | Galactokinase | ORFB01874 | galactokinase |
| ORF02112 | galactose operon repressor | NTORFA1667 | GalR, member of GalR-LacI family of transcriptional regulators, binds DNA; regulator of gal operon | ORFB01875 | sugar-binding transcriptional regulator, LacI family |
| ORF02113 | alcohol dehydrogenase, zinc-containing | NTORFA1668 | Alcohol dehydrogenase | ORFB01876 | alcohol dehydrogenase, zinc-containing, putative |
| | | | | ORFB01877 | alcohol dehydrogenase, zinc-containing, putative |

XY plots

Conserved synteny across multiple genomes.
Color highlighting of rearrangements at
whole genome level of resolution

# SNP report for reference sequence gmt_3810_assembly

**Please note:** if "(-)" appears in the query position column, then the query base reported has been reverse complemented.

| Ref asmbl | Ref pos | Ref type | Ref coding info | Query asmbl | Query pos | Query type | Query coding info | S/N |
|---|---|---|---|---|---|---|---|---|
| gmt_3810_assembly | 467 A | coding | ORF04243 MT0001 :: H (2) | ntmb01_2_assembly | 467 G | coding | NTORF0001 NTL01MB0001 :: R (2) | N |
| gmt_3810_assembly | 1057 G | coding | ORF04243 MT0001 :: V (1) | ntmb01_2_assembly | 1057 A | coding | NTORF0001 NTL01MB0001 :: I (1) | N |
| gmt_3810_assembly | 1849 C | intergenic | | bmt_689_assembly | 11067 A | intergenic | | |
| gmt_3810_assembly | 1977 G | intergenic | | ntmt02_1_assembly | 1977 A | intergenic | | |
| gmt_3810_assembly | 2347 A | coding | ORF04245 MT0002 :: D (2) | ntmb01_2_assembly | 2347 G | coding | NTORF0002 NTL01MB0002 :: G (2) | N |
| gmt_3810_assembly | 2532 T | coding | ORF04245 MT0002 :: L (1) | ntmb01_2_assembly | 2532 C | coding | NTORF0002 NTL01MB0002 :: L (1) | S |
| gmt_3810_assembly | 3751 T | coding | ORF04246 MT0003 :: L (1) | ntmb01_2_assembly | 3751 G | coding | NTORF0003 NTL01MB0003 :: V (1) | N |
| gmt_3810_assembly | 4013 C | coding | ORF04246 MT0003 :: T (2) | ntmt02_1_assembly | 4013 T | coding | NTORF0003 NTL02MT00003 :: I (2) | N |
| gmt_3810_assembly | 4480 C | coding | ORF04248 MT0004 :: S (2) | ntmb01_2_assembly | 4480 T | coding | NTORF0004 NTL01MB0004 :: L (2) | N |
| gmt_3810_assembly | 5752 G | coding | ORF04249 MT0005 :: V (3) | ntmb01_2_assembly | 5752 A | coding | NTORF0005 NTL01MB0005 :: V (3) | S |
| gmt_3810_assembly | 6406 C | coding | ORF04249 MT0005 :: N (3) | ntmb01_2_assembly | 6406 T | coding | NTORF0005 NTL01MB0005 :: N (3) | S |
| gmt_3810_assembly | 6446 G | coding | ORF04249 MT0005 :: A (1) | ntmb01_2_assembly | 6446 T | coding | NTORF0005 NTL01MB0005 :: S (1) | N |
| gmt_3810_assembly | 7362 C | coding | ORF04251 MT0006 :: Q (1) | ntmt02_1_assembly | 7362 G | coding | NTORF0006 NTL02MT00006 :: E (1) | N |
| gmt_3810_assembly | 7585 C | coding | ORF04251 MT0006 :: T (2) | ntmt02_1_assembly | 7585 G | coding | NTORF0006 NTL02MT00006 :: S (2) | N |
| gmt_3810_assembly | 8285 C | coding | ORF04251 MT0006 :: I (3) | ntmb01_2_assembly | 8285 T | coding | NTORF0006 NTL01MB0006 :: I (3) | S |

The selected SNP appears below the "SNP" indicator and is highlighted in red. Other sequence variations (which may or may not correspond to annotated SNPs) are also highlighted in red. Any ORFs or exons that overlap the displayed region are shown in blue, alongside the appropriate frame of the 6-frame conceptual protein translation. Use the link at the bottom of the page to expand the display to show the overlapping exons/ORFs in their entirety.

```
                                      S
                                      N
                                      P
                                      |
                                      |

    +1 N  V  G  Y  P  Q  S  G  R  H  P  C  G  R  A  P
    +2 M  W  D  I  R  N  R  G  V  I  P  A  G  A  L  P
    +2 M  W  D  I  R  N  R  G  V  I  P  A  G  A  L  P      gmt.ORF04280_cds;MT0026/+ strand
    +3  C  G  I  S  A  I  G  A  S  S  L  R  A  R  S
 28235 AATGTGGGATATCCGCAATCGGGGCGTCATCCCTGCGGGCGCGCTCCCC 28284    gmt_3810_assembly/+ strand
    -1  H  P  I  D  A  I  P  A  D  D  R  R  A  R  E  G
    -2  I  H  S  I  R  L  R  P  T  M  G  A  P  A  S  G
    -3   T  P  Y  G  C  D  P  R  *  G  Q  P  R  A  G


    +1 N  V  G  Y  P  Q  S  G  C  H  P  C  G  R  A  P
    +2 M  W  D  I  R  N  R  G  V  I  P  A  G  A  L  P
    +2 M  W  D  I  R  N  R  G  V  I  P  A  G  A  L  P      bmt.ORF00038_cds;/+ strand
    +3  C  G  I  S  A  I  G  V  S  S  L  R  A  R  S
 10886 AATGTGGGATATCCGCAATCGGGGTGTCATCCCTGCGGGCGCGCTCCCC 10935    bmt_690_assembly/+ strand
    -1  H  P  I  D  A  I  P  T  D  D  R  R  A  R  E  G
    -2  I  H  S  I  R  L  R  P  T  M  G  A  P  A  S  G
    -3   T  P  Y  G  C  D  P  H  *  G  Q  P  R  A  G


       AATGTGGGATATCCGCAATCGGGGTGTCATCCCTGCGGGCGCGCTCCCC        alternative_base [gray = unanimous bas

       00000000000000000000000000000000000000000000000000       read_depth [bmt_690_assembly min=1 max
       88888888888888888888888888888888888888888888888888

       22222222222223222323333322222222222223332222232333       consensus_quality_value [bmt_690_assem
       57787898643671688091343265565558676781019979090322

                                      |
                                      |
                                      S
                                      N
                                      P
```

# Sybil Completion

- Data model – 100%
- Workflow process – 100%
- Loading systems – 100%
- Interfaces – 80%, on-going
- SNP analysis – early phase
- Higher-level analysis systems – early phase

# Automated correction of genome sequence errors

Pawel Gajer[*], Michael Schatz and Steven L. Salzberg

The Institute for Genomic Research, 9712 Medical Center Drive, Rockville, MD 20850, USA

*To whom correspondence should be addressed. Tel: +1 301 795 7854; Fax: +1 301 795 7208; Email: pgajer@tigr.org

By using information from an assembly of a genome, a new program called AutoEditor significantly improves base calling accuracy over that achieved by previous algorithms. This in turn improves the overall accuracy of genome sequences and facilitates the use of these sequences for polymorphism discovery. We describe the algorithm and its application in a large set of recent genome sequencing projects. The number of erroneous base calls in these projects was reduced by 80%. In an analysis of over one million corrections, we found that AutoEditor made just one error per 8828 corrections. By substantially increasing the accuracy of base calling, AutoEditor can dramatically accelerate the process of finishing genomes, which involves closing all gaps and ensuring minimum quality standards for the final sequence. It also greatly improves our ability to discover single nucleotide polymorphisms (SNPs) between closely related strains and isolates of the same species.

# SNP Detection Pipeline Output

```
>KRUGERB-16 (94989bp) 200731 c-->t 50878 COV: 16 CB_QVal: 504 QUAL: 31.5
NON-SYNONYMOUS (2): GGC - G ---> GAC - D    (from start: 4 bp/848 bp)
REF : GAAGAAGCGCTAATAAAATGCCCATTACAAGACTCCCTTCG 200751 B. anthracis Ames Porton
        ||||||||||||||||||||| ||||||||||||||||||||
ASM : GAAGAAGCGCTAATAAAATGTCCATTACAAGACTCCCTTCG 50898
ORF01519          transporter  putative
  50873    A        564      AAAAAAAAAAAAAAAA    35:36:35:35:36:36:36:36:35:27:36:36:36:36:37:36
  50874    A        569      AAAAAAAAAAAAAAAA    34:36:35:35:36:36:36:36:35:35:36:36:36:36:35:36
  50875    A        550      AAAAAAAAAAAAAAAA    18:36:35:35:36:36:36:36:34:36:36:35:36:34:35:36
  50876    T        545      TTTTTTTTTTTTTTTT    13:36:36:35:36:36:36:34:35:38:35:34:36:34:37:34
  50877    G        517      GGGGGGGGGGGGGGGG    17:35:35:35:34:37:36:36:32:26:34:34:33:36:22:35
 *50878    T        504      TTTTTTTTTTTTTTTT    11:36:35:33:34:36:31:34:35:19:31:35:31:36:33:34
  50879    C        559      CCCCCCCCCCCCCCCC    33:34:35:35:36:36:34:35:32:35:36:35:36:36:35:36
  50880    C        515      CCCCCCCCCCCCCCCC    34:35:35:36:34:36:36:29:34:36:35:35:36:35:29
  50881    A        508      AAAAAAAAAAAAAAAA    36:34:34:36:40:34:36:32:26:36:34:36:36:32:26
  50882    T        513      TTTTTTTTTTTTTTTT    36:35:35:36:35:36:36:36:25:36:37:35:36:29:30
  50883    T        498      TTTTTTTTTTTTTTTT    36:35:34:36:33:36:36:35:15:36:35:32:36:33:30
REFERENCE: B. anthracis Ames Porton
  200726           A        334      AAAAAAAAAA      29:29:37:31:51:22:45:40:35:15
  200727           A        362      AAAAAAAAAA      25:41:37:34:51:27:45:40:34:28
  200728           A        347      AAAAAAAAAA      29:41:37:34:45:27:45:40:34:15
  200729           T        341      TTTTTTTTTT      25:41:45:33:45:21:45:25:34:27
  200730           G        336      GGGGGGGGGG      33:41:38:33:45:16:45:24:34:27
 *200731          C        334      CCCCCCCCCC      33:37:40:37:45:16:18:45:35:28
  200732           C        318      CCCCCCCCCC      33:37:40:38:45:18:4:45:35:23
  200733           C        284      CCCCCCCCC       33:37:37:37:41:15:4:45:35
  200734           A        307      AAAAAAAAA       27:45:37:37:51:21:4:51:34
  200735           T        291      TTTTTTTTT       24:45:37:37:41:24:4:45:34
  200736           T        293      TTTTTTATT       24:45:45:37:41:24:4:45:35
COV: 10 CBQUAL: 334 QUAL: 33.4
```

# Automated correction of genome sequence errors

**Pawel Gajer**[*], **Michael Schatz** and **Steven L. Salzberg**

The Institute for Genomic Research, 9712 Medical Center Drive, Rockville, MD 20850, USA

*To whom correspondence should be addressed. Tel: +1 301 795 7854; Fax: +1 301 795 7208; Email: pgajer@tigr.org

- ▸ **Full Text of this Article**
- ▸ Reprint (PDF) Version of this Article
- ▸ Email this article to a friend
- ▸ Similar articles found in:
  Nucl. Acids. Res. Online
  PubMed
- ▸ PubMed Citation
- ▸ Search PubMed for articles by:
  Gajer, P. || Salzberg, S. L.
- ▸ Alert me when:
  new articles cite this article
- ▸ Download to Citation Manager

By using information from an assembly of a genome, a new program called AutoEditor significantly improves base calling accuracy over that achieved by previous algorithms. This in turn improves the overall accuracy of genome sequences and facilitates the use of these sequences for polymorphism discovery. We describe the algorithm and its application in a large set of recent genome sequencing projects. The number of erroneous base calls in these projects was reduced by 80%. In an analysis of over one million corrections, we found that AutoEditor made just one error per 8828 corrections. By substantially increasing the accuracy of base calling, AutoEditor can dramatically accelerate the process of finishing genomes, which involves closing all gaps and ensuring minimum quality standards for the final sequence. It also greatly improves our ability to discover single nucleotide polymorphisms (SNPs) between closely related strains and isolates of the same species.